

Remarks

Regarding point 2 of the previous Office Action Previously filed claim 26 was rejected under 37 CFR 1.75(c) as being an improper dependent claim for failing to further limit the subject matter of a previous claim. The limitation "PCR primers" was rejected as an intended use limitation in a product claim. Applicants have responded by amending claim 26. The claim no longer recites the limitation "PCR primers".

Regarding point 3 of the previous Office Action

Claims 29, 31, 43 and 51 were rejected as indefinite under point 3 (A) due to the language "such as". Applicants have responded by amending these claims and removing the language "such as". Further discussion of this language is given below.

Claims 29, 31, 43 and 51 were further rejected as indefinite under point 3 (B) because of the language "paleospecies", "ecospecies", "agamospecies" and "ecosystem species". Clarification was requested. The applicants offer the following clarification. The applicants respectfully submit that the scope of these terms is definite. Paragraphs [0057] and [0058] of the specification relate to species and creatures. The terms appear in paragraph [0058]. These terms are known in the art of genetics. Applicants respectfully refer the Examiner to the enclosed definition of the terms "hybrid" and "species" on pages 188 and 364 respectively of the book A Dictionary of Genetics, 3rd Edition, (Oxford University Press, 1985, eds. R.C. King and W.D. Stansfield). The species definition uses the language above and further defines the language. Paragraph [0057] of the specification states *"The term creature means any organism that is living or was alive at one time."*

Page 188 of A Dictionary of Genetics defines the term "hybrid" as "an offspring from genetically dissimilar parents, perhaps even different species". Paragraph [0058] of the specification makes clear that the term "species hybrid" refers specifically to an offspring of different species, such as mules. Applicants respectfully submit that the language "species hybrid such as mules" is definite in the art of genetics, see p. 57 ("Animal Species Hybrid") and p. 469 (the "Hinny" or "Mule") of the Genetics Manual (G. P. Redei, World Scientific, 1998). In order to expedite claim allowance, however, the applicants have removed the language "such as mules" from the claims. The term "species hybrid" also includes plant species hybrids. See pp. 359-361 of Hybrid Origins of Plant Species (Annu. Rev. Ecol. Syst. 1997, 28: 359-89 by Loren Rieseberg); see specifically "What is a Hybrid Species?" Bottom p. 360, top p. 361.

Claims 44-51, especially independent claim 44 were rejected as indefinite under point 3 (C) due to the language *"is used by the apparatus to determine the data"*. The applicants have amended the claim by deleting this phrase as was suggested by the Examiner in the previous Office Action.

Some further remarks

The applicants hereby respectfully submit some further remarks to aid the examination of the claims.

Regarding claim 7 the limitation *"and a population and the population is a group of individuals as in the field of population genetics"* has been eliminated from this claim and other claims. As seen from paragraphs [0175]-[0176], a CL-F region need not necessarily be limited by such a limitation.

Regarding claim 8 the limitation *"wherein the width of the subrange of the segment-subrange is less than 0.5 and whereby the segment of the segment-subrange is a chromosome segment and the length of the chromosome segment is less than or equal to the length of the chromosome"* has been added to this claim and others. As seen from paragraphs [0062], [0063] and [0095], the specification necessarily describes subranges within and smaller than the range 0 to 0.5. The widths of these subranges are necessarily less than 0.5, i.e. $0.5 = 0.5 - 0$. The length of a chromosome segment is less than or equal to the length of a chromosome, see [0275], bottom p. 20.

Regarding claim 17 the amended claim recites *"A copy of a set of oligonucleotides.... wherein each oligonucleotide in the set is a type (1) complementary oligonucleotide or wherein each oligonucleotide in the set is a type (2) complementary oligonucleotide, wherein each bi-allelic covering marker is an exact, true bi-allelic marker"*. Applicants respectfully submit that "a copy" is described by "one or more copies" [0255], [0265]. Type (1) and type (2) complementary oligonucleotides are described in paragraphs [0142] and [0143].

Type (1) oligonucleotides are used, for example, as sequence specific type oligonucleotides. Examples of their use are given in [0324] and [0349]. Type (2) oligonucleotides are used, for example, as standard PCR primers see [0143]. Some examples of their use is given in [0144] and [0249].

The term "bi-allelic markers" in the art generally means exact, true bi-allelic markers. (For example, SNPs are examples of such exact, true bi-allelic markers.) The specification, however, also expands the term "bi-allelic marker" somewhat and describes bi-allelic marker equivalents or BMEs (mathematical markers formed from one or more markers that act like they are bi-allelic) and approximate bi-allelic markers, see for example [0054] and [0055].

Regarding claim 26

This claim includes the newly added limitation “*wherein thousands of the covering markers are from one chromosome*”. Other claims such as claims 38, 49, and 52 also include this limitation or a similar limitation.

As stated on page 18 of the previously filed Supplemental Amendment/Response of Nov. 20, 2005, at the time the application was filed the whole field of association studies is looking to use thousands of bi-allelic markers. See for example Risch, N. and Merikangas, K.: The Future of Genetic Studies of Complex Human Diseases. Science, 13 September 1996, vol. 273, pp. 1516-1517 cited in [0027] of the application. This Risch paper (see p. 1517 mid left most column) describes using technological advances to do association testing of five diallelic (or bi-allelic) polymorphisms within each of 100, 000 genes (a total of 500, 000 polymorphisms tested in the association study). **This is 20, 000 or more markers on a chromosome.** And the inventor's paper is a generalization of the Risch and Merikangas analysis [0029]. A copy of the Risch paper was supplied with the Amend/Resp of 11/05 and is also included with this document.

Another example of the expectation (at the time of filing) in the field of using large numbers of markers (e.g., thousands) from high-density marker maps is the Kruglyak paper. The Kruglyak paper (*The use of a genetic map of bi-allelic markers in linkage studies*, published 9/97, see footnote 4, p.3 of the present application) is quoted in the application (see mid [0026]) as predicting a density of at least 1,000 SNPs (bi-allelic markers) per cM. **Since a human chromosome is about 150 cM in length, a density of at least 1,000 bi-allelic markers per cM is about at least 150, 000 bi-allelic markers on a chromosome (or at least 3 million bi-allelic markers in the genome).** Page 21 (right column) of the Kruglyak paper essentially predicts that this large number (and density) of markers could be practically genotyped using more automated genotyping techniques (p. 24 Kruglyak, endnotes 10, 12, 14, 15, 16) that are also described in the present application under Oligonucleotide Technology [0249] in endnote 11 (application p. 25) in references (1) Chee, (2) Saiki, (3) Wu, and (4) Nickerson. (A copy of the Kruglyak paper was supplied with the Amend/Resp of 11/05 and is also included with this document).

Another example of this expectation in the field is the Chee paper. The Chee paper ([0341], endnote 8) is part of this application through incorporation by reference. The Chee paper describes a high-density array (or "gene chip") *"that could query the entire coding content of the human genome, estimated at 100, 000 genes"* (see p. 613, last sentence of the paper). **A total of 100, 000 genes necessarily means more than 5, 000 markers per chromosome.** The Chee paper also describes high-resolution marker maps (p. 613 left most column). The Chee paper is cited in the application in connection with using thousands of bi-allelic markers in the new two-dimensional techniques of this application, see [0249] and [0324]. More specifically paragraph [0325] describes the use of "gene chips" as a physical implementation used to scan a particular chromosome or chromosomal region. Since the markers for scanning the particular chromosome or chromosomal region must be from the chromosome, the limitation *"wherein thousands of the covering markers are from one chromosome"* is supported. (A copy of pages of the Chee paper (pp. 610 and 613) were included for the Examiner's convenience in the previously filed Amendment/Response of 9/13/05 and are also included with this document.)

In addition, the specification describes using more dense marker coverings in [0182] and [0183], wherein N is large and the covering distance δ is small. And these are described in conjunction with larger CL-F regions that are N covered. Examples of CL-F regions that are N covered include segment-subranges [0185] (wherein, as stated above, the segment covered is less than or equal in length to the length of a chromosome). And other examples are CL-F regions wherein the chromosomal location coordinates of the CL-F region range over a chromosome or part of a chromosome [075] further support the limitation *"wherein thousands of the covering markers are from one chromosome"*.

Regarding claim 27 the limitations 0.2 and 12 cM are supported for example by [0180] and [0181].

Regarding claim 37 this claim is the same scope as previously allowed old claim 38. As is apparent from the claim listing, claim 37 simply incorporates the limitation *"wherein $N > 2$ "* from previously allowed old claim 38.

Regarding claim 46 the limitation *"oligonucleotides bound to a glass slide of silicon chip"* is already present as one limitation in claim 44, from which claim 46 depends. See also paragraphs [144] and [0323]. The limitation *"wherein each covering marker is an exact, true bi-allelic"* marker is discussed above.

Conclusion

This Amendment/Response has responded to each point of rejection in the previous Office Action of 12/29/05. The applicants have amended several claims and five new claims have also been added. Appropriate fees are also enclosed.

For the reasons advanced above, applicants respectfully submit that the application is now in condition for allowance and that action is earnestly solicited.

Respectfully submitted,



Robert O. McGinnis
Registration No. 44, 232

May 30, 2006
1575 West Kagy Blvd.
Bozeman, MT. 59715
tel (406)-522-9355

Enclosures:

- 1) A Dictionary of Genetics (3rd Edition, Oxford University Press, 1985) Title page and pp. 188 and 364. (3 sheets total)
 - 2) Genetics Manual (World Scientific, 1998) Title page and pp. 57 and 469. (3 sheets total)
 - 3) Hybrid Origin of Plant Species by Loren Rieseberg (Annu. Rev. Ecol. System. 1997. 28:359-89) pp. 359-361 (3 sheets total)
 - 4) Future of Genetic Studies of Complex Human Diseases by Risch, N. & Merikangas, K. (Science vol. 273 13 September 1996 pp. 1516-1517) pp. 1516 & 1517 (total 2 sheets)
 - 5) The use of a genetic map of biallelic markers in linkage studies by L. Kruglyak (Nature Genetics vol. 17, Sept. 1997, pp. 21-24) pp. 21-24 (total of 4 sheets)
 - 6) Accessing Genetic Information with High-Density DNA Arrays by Chee, M. et. al., (Science vol. 274, 25 Oct. 1996, pp. 610-614) pp. 610 and 613 (2 sheets)
- 17 total enclosure sheets

A **DICTIONARY** **OF GENETICS**

THIRD EDITION

ROBERT C. KING

PROFESSOR OF BIOLOGY
NORTHWESTERN UNIVERSITY

WILLIAM D. STANSFIELD

PROFESSOR OF BIOLOGY
CALIFORNIA POLYTECHNIC STATE UNIVERSITY

New York Oxford
OXFORD UNIVERSITY PRESS
1985

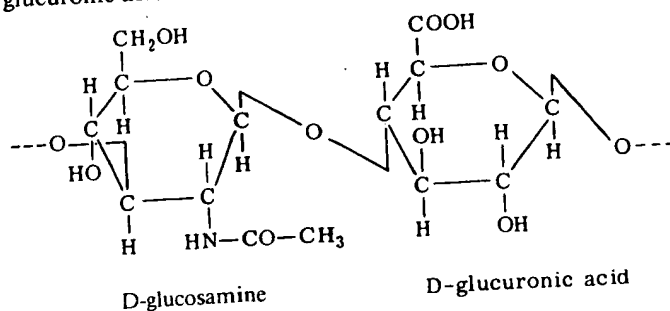
BEST AVAILABLE COPY

H-X antigen a histocompatibility antigen encoded by a locus on the X chromosome.

Hyalophora cecropia the giant cecropia moth; because of its large size a favorite experimental insect.

hyaloplasm cytosol.

hyaluronic acid a mucopolysaccharide that is abundant in the jelly coats of eggs and in the ground substance of connective tissue. Hyaluronic acid is a polymer composed of glucosamine and glucuronic acid subunits.



hyaluronidase an enzyme that digests hyaluronic acid.

H-Y antigen an antigen detected by cell-mediated and humoral responses of homogametic individuals against heterogametic individuals of the same species, which are otherwise genetically identical. Antigenic responses of this sort have been demonstrated in mammals, birds, and amphibians. In mammals the antigen is called H-Y because it acts as a *Histocompatibility* factor determined by the Y-chromosome. The H-Y antigen is the major male-determining factor in the developing mammalian male fetus. The location of the gene encoding the H-Y antigen is not known. However, the gene which induces synthesis of the H-Y antigen in humans is located on the short arm of the Y-chromosome near the centromere. A homologous locus, which suppresses H-Y production, lies on the distal end of the short arm of the X. Evidently two doses of this gene are necessary for the complete suppression, since Turner syndrome (45, XO) females produce small amounts of the H-Y antigen. The H-Y locus is one of the areas that escapes X-chromosome inactivation (*q.v.*).

* **hybrid** 1. a heterozygote (e.g., a monohybrid is heterozygous at a single locus; a dihybrid is heterozygous at two loci; etc.). 2. an offspring from genetically dissimilar parents, perhaps even different species.

hybrid arrested translation a method for identifying the cDNA corresponding to an mRNA that depends upon the ability of cDNA to hybridize with its mRNA and thereby to inhibit its translation in an *in vitro* system; the disappearance of the translation product from the system indicates the presence of the cDNA.

hybrid breakdown the reduction in fitness of F_2 and/or backcross populations from fertile hybrids produced by intercrossing genetically disparate populations or species; a postzygotic reproductive isolating mechanism.

hybrid corn commercial corn grown from seed produced by the "double cross" (*q.v.*) procedure. Such corn is characterized by its vigor and uniformity.

hybrid DNA model a model used to explain both crossing-over and gene conversion by postulating that a short segment of heteroduplex (hybrid) DNA is produced from both parental DNAs in the neighborhood of a chiasma. See *Holliday model*.

hybrid duplex molecule an experimentally reconstituted molecule containing a segment

of single-stranded DNA, complementary base sequence.

hybrid dysgenesis a phenomenon which occurs spontaneously when hybrids of different germ line defects including mutations, and in most cases possible element name strains susceptible to strains established from dysgenic F_1 individual P strains their transposons placed in M cytoplasmic background at high rates, dihybrid element.

hybrid inviability a phenomenon observed in hybrids between dissimilar species.

hybridization 1. the process of hybridizing two or to different species. 2. the pairing of phenotypes. 3. the pairing of DNA hybrid, or the sources.

hybridization competent a cell or tissue capable of using a variation of the technique of trapping on a nitrocellulose membrane to that DNA. An unlabeled DNA, if it will complement the diminution of labeling by the unlabeled DNA.

hybridoma a cell (usually a lymphocyte) and a myeloma cell clone that can be made to secrete only a single type of antibody.

hybrid resistance the resistance of hybrid recipients than in the parental histocompatible with.

hybrid sterility the sterility of the offspring.

hybrid swarm a collection of hybridization of two or more generations.

hybrid tobacco mosaic a disease caused by acid and protein compounds.

hybrid vigor heterosis.

hybrid zone a geographical area where two species are observed.

hydrocarbon an organic compound consisting of carbon and hydrogen atoms.

SOS boxes the operator sequences in *E. coli* DNA that are recognized by a repressor called the LexA protein. This protein represses several loci involved in DNA repair functions. See **SOS response**.

SOS response an error-prone mechanism of repairing damaged DNA in *E. coli* by the coordinated induction of several enzymes. Damaged DNA somehow activates an enzyme called RecA protease, and this protease cleaves a protein called LexA repressor. Many genes involved in repair functions become activated when this repressor is cleaved. See **SOS boxes**.

Southern blotting a technique, developed by E.M. Southern, for transferring electrophoretically resolved DNA segments from an agarose gel to a nitrocellulose filter paper sheet via capillary action. Subsequently the DNA segment of interest is probed with a radioactive, complementary nucleic acid, and its position is determined by autoradiography. A similar technique, referred to as *northern blotting*, is used to identify RNAs. For example, an electropherogram containing a multitude of different mRNAs could be probed with a radioactive cloned gene. In cases where proteins have been separated electrophoretically, a specific protein on an electropherogram can be identified by the *western blotting* procedure. In this case the probe is a radioactively labeled antibody raised against the protein in question.

sow the adult female of swine.

spacer DNA untranscribed segments of eukaryotic and some viral genomes flanking functional genetic regions (cistrons). Spacer segments usually contain repetitive DNA. The function of spacer DNA is not presently known, but it may be important for synapsis. See **transcribed spacer**.

spawn to deposit eggs.

spay to remove the ovaries.

special creation a nonscientific philosophy asserting that each species has originated through a separate act of divine creation by processes that are not now in operation in the natural world.

specialized 1. an organism having a narrow range of tolerance for one or more ecological conditions. 2. a species having a relatively low potential for further evolutionary change; the opposite of generalized.

specialized transduction See **transduction**.

speciation 1. the splitting of an ancestral species into daughter species that co-exist in time; horizontal evolution or speciation; cladogenesis. 2. the gradual transformation of one species into another without an increase in species number at any time within the lineage; vertical evolution or speciation; phyletic evolution or speciation.

* **species** 1. biological (genetic) species: reproductively isolated systems of breeding populations. 2. paleospecies (successional species): distinctly different appearing assemblages of organisms as a consequence of species transformation (*q.v.*). 3. taxonomic (morphological; phenetic) species: phenotypically distinctive groups of coexisting organisms. 4. microspecies (agamospecies): asexually reproducing organisms (mainly bacteria) sharing a common morphology and physiology (biochemistry). 5. biosystematic species (ecospecies; coenospecies): populations that are isolated by ecological factors rather than ethological isolation (*q.v.*).

species group superspecies (*q.v.*).

species selection a form of group selection (*q.v.*) in which certain species (produced by cladogenesis) continue the cladogenic process and others become extinct.

species
the pas
species
specific
same k
stable e
specific
ticular
exposur
specific
in a give
specifici
strate, b
respondi
specificit
polymer
(*q.v.*). See
specimen
sisting of
spectroph
of light
medium.
spelt *Tr*
the latter
S period
sperm a
spermatel
spermathe
matzoa d
spermatid
out furthe
(*q.v.*).
spermatocy
spermatocy
ary spermat
spermatogei
spermatogoi
genitors of
Spermatoph
contempora
len tubes and
See classific
spermatozoa
sperm bank

BEST AVAILABLE COPY

GENETICS MANUAL

CURRENT THEORY, CONCEPTS, TERMS

GEORGE P. RÉDEI

University of Missouri



World Scientific

Singapore • New Jersey • London • Hong Kong

RÉDEI

Animal models continued

pigmentation of the retina, h. chr. 6p21.2-cen, mouse gene *RD2^{Rd2}*, m. chr. 17), *gonadal dysgenesis* (underdeveloped germcells in the testes, h. chr. Y11.2-pter, mouse gene *Sry^{Sxr}*, m. chr. Y), *tyrosinase negative oculocutaneous albinism* (see albinisms, h. chr. 11q14-q21, mouse gene *Tyrc*, m. chr. 7). By disruption of hexosaminidase α subunit a model for the Tay-Sachs disease has been generated in mouse. Interestingly, these animals suffered no obvious behavioral or neurological deficit. Disrupting the hexosaminidase β subunit (Sandhoff disease model) resulted in massive depletion of spinal cord axons and neuronal storage of ganglioside G_{M2} .

MOUSE POLYGENIC DISORDERS WITH SIMILARITIES TO HUMAN CONDITIONS [human problem - mouse strain]: alcoholism and opiate drug addictions - C57BL/6J, asthma - A/J, atherosclerosis - C57BL, audiogenic (sound-induced) seizures - DBA, cleft palate (fissure in the mouth) - A, deafness - LP, dental disease - C57BL, BALB/c, diabetes - NOD, epilepsy - EL, SWXL-4, granulosa cell tumors in the ovary - SWR, germ cell tumors in the ovary - LT, germ cell tumors in the testes - 129, hemolytic anemia - NZB, hepatitis - BALB/c, Hodgkin disease (pre-B cell lymphoma - SJL, hypertension - MA/My, kidney adenocarcinoma - BALB/cCd, leprosy (*Mycobacterium leprae*) - BALB/c, leukemia - AKR/J, C58/J, P/J, lung tumors - A, Ma/My, measles - BALB/c, osteoporosis - DBA, polygenic obesity - NZB, NZW, pulmonary tumors - A/J, rheumatoid arthritis - MRL/Mp, spina bifida (defect of the bones of the spinal cord) - CT, systemic lupus erythematosus (a skin degeneration) - NZB, NZW, whooping cough (pertussis) - BALB/c. (Some of the data by courtesy of GIBCO BRL Co.)

ANIMAL POLE: is dorsal end of the animal egg opposite the lower end, the vegetal pole, and where the sperm entry is located. After the entry, the egg cortex rotates slightly and in some species at the side opposite the entry a *gray crescent* is formed. (See vegetal pole)

ANIMAL SPECIES HYBRIDS: the most familiar example is the hybrids of the mare (*Equus caballus*, $2n = 64$) and the jackass (*Equus asinus*, $2n = 62$), and the stallion and the she-ass. The hybrid males do not produce viable sperm although they may show normal libido. The females may have estrus and ovulate but there is no proven cases of fertility. Zebras ($2n = 44$) also may form hybrids with both donkeys and horses. Buffalo (*Bison bison*, $2n = 60$) may be crossed reciprocally with cattle (*Bos taurus*, $2n = 60$) but their offspring (cattalo) has reduced fertility. The domesticated pig (*Sus crofa*, $2n = 38$) forms fertile hybrids with several wild pigs



with the same number of chromosomes. The sheep (*Ovis aries*, $2n = 54$) interbreeds with the wild mouflons but the sheep x goat (*Capra hircus*, $2n = 60$) hybrid embryo only rarely can be kept alive. Some monkeys can be interbred but primates are generally sexually isolated. There is no sexual barrier among the various human races, indicating close relationship but no hybrids are known between humans and any other species. These general rules do not hold for somatic cell hybrids because human cells can be fused with rodent or plant cells but they cannot be regenerated or even maintained successfully for indefinite periods of time. The hybridization barrier is not identical with other functional barriers.

HYBRID OF THE MALE GRANT'S ZEBRA AND THE FEMALE BLACK ARABIAN ASS, GLOUCESTER ZOO. (From Gray, A.P. Mammalian Hybrids. Commonwealth Agric Bureau, Farnham Road, Slough, UK)

ANIMAL TRANSFORMATION VECTORS: most commonly Simian virus 40 (SV40) and Bovine papilloma virus (BPV) based vectors are used. The BPV vectors can be used for the

RÉDEI

High-lysine corn continued

tein the lysine contents are ca. : gliadin 5.0, glutenin 17.6, albumin 78.4 and globulin 98.0. Cereals with improved nutritional values are desirable for feeding of animals and more importantly for the production of cereal food for human populations suffering from malnutrition as a result of protein deficiency in the diet (kwashiorkor). (See also essential amino acids, kwashiorkor)

HIGH-MOBILITY GROUP OF PROTEINS (HMG): are associated with functionally active chromatin and render the genes more sensitive to DNase and probably to RNase II. They are regulated by cell-cycle-dependent phosphorylation, affecting their ability to bind to DNA. HMG proteins are important for growth and development. A large number of transcription factors contain HMG-like domains. The specificity of these proteins varies but a common feature is that they distort DNA structure and have an affinity for distorted DNA structures. The change in DNA electrophoretic mobility is correlated with this altered structure. Recurrent rearrangements of the HMG1-C group was detected in some benign tumors. (See chromatin, nonhistone proteins, cell cycle, coactivator, transcription factors, SOX, DNA bending)

HIGH-PERFORMANCE LIQUID CHROMATOGRAPHY (HPLC): a mixture of compounds are applied to chromatographic columns with strong ion exchange resins. The solvent is forced through the resin under pressure for rapid and sharp separation of the components of the mixture. The eluates are electronically scanned and identified. (See chromatography)

HIGHLY REPETITIVE DNA: contains high degree of redundancy and reassociates very rapidly after denaturation. (See SINE, LINE, annealing, c_0t value)

HILL REACTION: illuminated chloroplasts evolve oxygen and reduce an artificial electron acceptor (ferricyanide \rightarrow ferrocyanide). It was an important tool to study the mechanism of photosynthesis, namely that the evolved oxygen comes from water rather than from CO_2 and demonstrated that isolated chloroplasts can perform part of the reactions, and revealed the light-activated transfer of an electron from one substance to another against a chemical-potential gradient. (See photosynthesis)

HILUM: a depression or pit where vessels and nerves enter an organ; the place where the plant seed is connected to its stalk in the fruit.

him: high incidence of males mutation in *Caenorhabditis* have a high level of nondisjunction in XX hermaphrodites and thus produce XO males. (See nondisjunction, chromosomal sex determination, *Caenorhabditis*)

HIMALAYAN RABBIT: carries temperature-sensitive tyrosinase genes controlling pigmentation. The extremities, paws, ears and tail having lower blood circulation and concomitant lower body temperature develop darker pigmentation. Similar pattern of pigmentation occurs in other rodents and in the Siamese cats. Tyrosinase is a copper enzyme (also called polyphenol oxydase) is involved in the formation of 3,4-dihydroxyphenylalanine (DOPA) that is responsible for the production of melanin in the hair and skin and darkening of wounded fruits and other plant tissues. (See albinism, piebaldism, Siamese cat, pigmentation of animals, temperature-sensitive mutation)

HindIII: restriction enzyme with recognition site $A\downarrow AGCTT$.

HiNF (histone nuclear factor): a 48 K M_r protein, identical to interferon regulatory factor IRF-2. (See IRF-2, histones)

HINGE: see antibody

* **HINNY:** she-ass ($2n = 62$) x stallion ($2n = 64$) hybrid. The reciprocal (mare x jackass) is called mule. Mules are easier to produce because the jackass willingly mates with the mare but the stallion mates with the she-ass only under special circumstances (blindfolded). The hybrids' body resembles closer the female parent as an apparent cytoplasmic influence. These sterile hybrids—known since the beginning of human civilization—may retain some sexual drive and occasionally fertility has been reported in backcrosses with either the jackass or the stallion. The jackass backcrosses are entirely sterile but the backcrosses with stallions appear more nor-



HIMALAYAN RABBIT

HYBRID ORIGINS OF PLANT SPECIES

Loren H. Rieseberg

Biology Department, Indiana University, Bloomington, Indiana 47405;
e-mail: lriesebe@bio.indiana.edu

KEY WORDS: plants, hybridization, introgression, reproductive isolation, speciation

ABSTRACT

The origin of new homoploid species via hybridization is theoretically difficult because it requires the development of reproductive isolation in sympatry. Nonetheless, this mode is often and carelessly used by botanists to account for the formation of species that are morphologically intermediate with respect to related congeners. Here, I review experimental, theoretical, and empirical studies of homoploid hybrid speciation to evaluate the feasibility, tempo, and frequency of this mode. Theoretical models, simulation studies, and experimental syntheses of stabilized hybrid neospecies indicate that it is feasible, although evolutionary conditions are stringent. Hybrid speciation appears to be promoted by rapid chromosomal evolution and the availability of a suitable hybrid habitat. A selfing breeding system may enhance establishment of hybrid species, but this advantage appears to be counterbalanced by lower rates of natural hybridization among selfing taxa. Simulation studies and crossing experiments also suggest that hybrid speciation can be rapid—a prediction confirmed by the congruence observed between the genomes of early generation hybrids and ancient hybrid species. The frequency of this mode is less clear. Only eight natural examples in plants have been rigorously documented, suggesting that it may be rare. However, hybridization rates are highest in small or peripheral populations, and hybridization may be important as a stimulus for the genetic or chromosomal reorganization envisioned in founder effect and saltational models of speciation.

INTRODUCTION

Hybridization may have several evolutionary consequences, including increased intraspecific genetic diversity (2), the origin and transfer of genetic adaptations

(2, 93), the origin of new ecotypes or species (42, 102), and the reinforcement or breakdown of reproductive barriers (27, 55, 77). Although the frequency and importance of these outcomes are not yet clear in either plants or animals, a critical body of data is now available for assessing the mechanistic basis and frequency of one of these—the origin of new species. The last comprehensive review of this topic in relation to plants was Grant's (42) monograph "Plant Speciation." Grant listed six mechanisms by which the breeding behavior of hybrids could be stabilized, thus providing the potential for speciation:

1. asexual reproduction;
2. permanent translocation heterozygosity;
3. permanent odd polyploidy;
4. allopolyploidy;
5. the stabilization of a rare hybrid segregate isolated by postmating barriers;
6. the stabilization of a rare hybrid segregate isolated by premating barriers.

The first three of these mechanisms generate flocks of clonal or uniparental microspecies that span the range of morphological variability between the parental species. Sexual reproduction among microspecies is limited or absent, making it difficult to discuss their origin and evolution in the context of sexual isolation and speciation. By contrast, the latter three mechanisms generate sexual derivatives and therefore have the potential to give rise to new biological species.

This review focuses on the origin of sexual, homoploid hybrid species (mechanisms 5 and 6), (but see 50, 89 for reviews of polyploidy in plants). After clarification of concepts and terminology, the historical basis of our current understanding of hybrid speciation is reviewed. This is followed by examination of the frequency of natural hybridization and an exploration of experimental and theoretical studies that test the feasibility of homoploid hybrid speciation. Once the feasibility of this mode of speciation has been established, I briefly critique the methods used for identifying homoploid hybrid species in nature and then focus on those examples of hybrid speciation that are well established. Finally, I discuss promising areas for future research and possible approaches that may facilitate studies of this mode.

✱ WHAT IS A HYBRID SPECIES?

Both "hybrid" and "species" can have several meanings for evolutionary biologists. The term hybrid can be restricted to organisms formed by cross-fertilization between individuals of different species, or it can be defined more

broadly as the offspring between individuals from populations "which are distinguishable on the basis of one or more heritable characters" (44). I prefer this broader definition of hybrids, as it provides greater flexibility in usage. Nonetheless, in this review, I focus on hybrids formed by crosses between species. ✱

✱ The term species has a much wider variety of definitions, ranging from concepts based on the ability to interbreed to those based on common descent. Mayr's (59) biological species concept—"species are groups of interbreeding natural populations which are reproductively isolated from all other such groups"—is perhaps the most widely accepted of these. Although I have previously expressed concern about the limitations of this concept (73), its emphasis on reproductive isolation does offer a straightforward approach to the study of speciation (20). Moreover, the evolution of reproductive barriers is particularly crucial to the successful origin of new hybrid species; otherwise, the new hybrid lineage will be swamped by gene flow with its parents. Thus, the focus of this review is on the evolution of reproductive isolation between new hybrid lineages and their parents.

HISTORICAL PERSPECTIVE

The hypothesis that new species may arise via hybridization appears to have originated with Linnaeus (58; cited in 84), who wrote "it is impossible to doubt that there are new species produced by hybrid generation. . . . For thence it appears to follow, that the many species of plants in the same genus in the beginning could not have been otherwise than one plant, and have arisen from this hybrid generation." This represents a modification of the orthodox view of special creation, which asserted that all existing species were created by the hand of God and which denied the existence of constant hybrids (15). However, Linnaeus' observations were limited to F_1 hybrids, and he was unaware of potential difficulties with his hypothesis such as segregation and sterility.

Rigorous experimental study of plant hybridization was initiated by Joseph Kölreuter in 1760 and led to two critical discoveries (84). First, Kölreuter found that a hybrid from *Nicotiana paniculata* \times *N. rustica* produced no seeds—the first "botanical mule." As a result, Kölreuter concluded that hybrid plants are produced only with difficulty and are unlikely to occur in nature in the absence of human intervention or disturbance to the habitat. Second, Kölreuter and his successor, Carl Gartner, discovered that later generation hybrids tended to revert back to the parental forms, thus refuting the existence of constant hybrids and supporting the orthodox view of special creation (84). The views of Kölreuter and Gartner on the lack of constancy of hybrids (although not necessarily on creation) were held by most other prominent botanical hybridizers during the eighteenth and nineteenth centuries, including Charles Darwin, John

The Future of Genetic Studies of Complex Human Diseases

Neil Risch and Kathleen Merikangas

Geneticists have made substantial progress in identifying the genetic basis of many human diseases, at least those with conspicuous determinants. These successes include Huntington's disease, Alzheimer's disease, and some forms of breast cancer. However, the detection of genetic factors for complex diseases—such as schizophrenia, bipolar disorder, and diabetes—has been far more complicated. There have been numerous reports of genes or loci that might underlie these disorders, but few of these findings have been replicated. The modest nature of the gene effects for these disorders likely explains the contradictory and inconclusive claims about their identification. Despite the small effects of such genes, the magnitude of their attributable risk (the proportion of people affected due to them) may be large because they are quite frequent in the population, making them of public health significance.

Has the genetic study of complex disorders reached its limits? The persistent lack of replicability of these reports of linkage between various loci and complex diseases might imply that it has. We argue below that the method that has been used successfully (linkage analysis) to find major genes has limited power to detect genes of modest effect, but that a different approach (association studies) that utilizes candidate genes has far greater power, even if one needs to test every gene in the genome. Thus, the future of the genetics of complex diseases is likely to require large-scale testing by association analysis.

How large does a gene effect need to be in order to be detectable by linkage analysis? We consider the following model: Suppose a disease susceptibility locus has two alleles A and a, with population frequencies p and $q = 1 - p$, respectively. There are three genotypes: AA, Aa, and aa. We define genotypic relative risks (GRR, the increased chance that an individual with a particular genotype has the disease) as follows: Let the risk for individuals of genotype Aa be γ times greater than the risk for individuals with genotype aa, a GRR of γ . We assume a multiplicative relation for two A alleles, so that the GRR for genotype AA is γ^2 . The method of link-

age analysis we have chosen for this argument is a popular current paradigm in which pairs of siblings, both with the disease, are examined for sharing of alleles at multiple sites in the genome defined by genetic markers. The more often the affected siblings share the same allele at a particular site, the more likely the site is close to the disease gene. Using the formulas in (1), we calculate the expected proportion Y of alleles shared by a pair of affected siblings for the best possible case—that is, a closely linked marker locus (recombination fraction $\theta = 0$) that is fully informative (heterozygosity = 1) (2)—as

$$Y = \frac{1 + w}{2 + w} \text{ where } w = \frac{pq(\gamma - 1)^2}{(p\gamma + q)^2}$$

If there is no linkage of a marker at a particular site to the disease, the siblings would be expected to share alleles 50% of the time; that is, Y would equal 0.5. Values of Y for various values of p and γ are given in the third column of the table. For an allele of moderate frequency (p is 0.1 to 0.5) that confers a GRR (γ) of fourfold or greater, there is a detectable deviation of Y from the null value of 0.5. On the other hand, for an allele conferring a GRR of 2 or less, the expected marker-sharing only marginally exceeds 50%, for any allele frequency (p). Thus, it is clear that the use of

linkage analysis for loci conferring GRR of about 2 or less will never allow identification because the number of families required (more than ~2500) is not practically achievable.

Although tests of linkage for genes of modest effect are of low power, as shown by the above example, direct tests of association with a disease locus itself can still be quite strong. To illustrate this point, we use the transmission/disequilibrium test of Spielman *et al.* (3). In this test, transmission of a particular allele at a locus from heterozygous parents to their affected offspring is examined. Under Mendelian inheritance, all alleles should have a 50% chance of being transmitted to the next generation. In contrast, if one of the alleles is associated with disease risk, it will be transmitted more often than 50% of the time.

For this approach, we do not need families with multiple affected siblings, but can focus just on single affected individuals and their parents. For the same model given above, we can calculate the proportion of heterozygous parents as $pq(\gamma + 1)/(p\gamma + q)$ (4). Similarly, the probability for a heterozygote parent to transmit the high risk A allele is just $\gamma/(1 + \gamma)$. Association tests can also be performed for pairs of affected siblings. When the locus is associated with disease, the transmission excess over 50% is the same as for single offspring, but the probability of parental heterozygosity is increased at low values of p ; for higher values of p , the probability of parental heterozygosity is decreased. The formula for parental heterozygosity for an affected pair of siblings for the same genetic model as used in the first example is

$$h = \frac{pq(\gamma + 1)^2}{2(p\gamma + q)^2 + pq(\gamma - 1)^2}$$

Linkage					Association			
Genotypic risk ratio (γ)	Frequency of disease allele A (p)	Probability of allele sharing (Y)	No. of families required (N)	Probability of transmitting disease allele A ($P(\text{tr-A})$)	Singletons		Sib pairs	
					Proportion of heterozygous parents (Het)	(N)	(Het)	(N)
4.0	0.01	0.520	4260	0.800	0.048	1098	0.112	235
	0.10	0.597	185	0.800	0.346	150	0.537	48
	0.50	0.576	297	0.800	0.500	103	0.424	61
	0.80	0.529	2013	0.800	0.235	222	0.163	161
2.0	0.01	0.502	296,710	0.667	0.029	5823	0.043	1970
	0.10	0.518	5382	0.667	0.245	695	0.323	264
	0.50	0.526	2498	0.667	0.500	340	0.474	180
	0.80	0.512	11,917	0.667	0.267	640	0.217	394
1.5	0.01	0.501	4,620,807	0.600	0.025	19,320	0.031	7776
	0.10	0.505	67,816	0.600	0.197	2218	0.253	941
	0.50	0.510	17,997	0.600	0.500	949	0.490	484
	0.80	0.505	67,816	0.600	0.286	1663	0.253	941

Comparison of linkage and association studies. Number of families needed for identification of a disease gene.

N. Risch is in the Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305-5120, USA. E-mail: risch@lahmed.stanford.edu. K. Merikangas is in the Departments of Epidemiology and Psychiatry, Unit, Yale University School of Medicine, New Haven, CT 06510, USA. E-mail: kath@zeus.psych.yale.edu

On the right side of the table, we present the proportion of heterozygous parents (Het) and the probability of transmission of the A allele from a heterozygous parent to an affected child $[P(\text{tr-A})]$ for the same values of GRR as considered above for the example of linkage analysis. The deviation from the null hypothesis of 50% transmission from heterozygous parents is substantially greater than the excess allele sharing that is found by linkage analysis in sibling pairs. This disparity between the methods is particularly true for lower values of γ (that is, with lower relative risk). For example, for $\gamma = 1.5$, allele sharing is at most 51%, while the A allele is transmitted 60% of the time from heterozygous parents.

In this respect then, association studies seem to be of greater power than linkage studies. But of course, the limitation of association studies is that the actual gene or genes involved in the disease must be tentatively identified before the test can be performed. In fact, the actual polymorphism within the gene (or at least a polymorphism in strong disequilibrium) must be available. However, we show that this requirement is only daunting because of limitations imposed by current technological capabilities, not because sufficient families with the disease are not available or the statistical power is inadequate (5). For example, imagine the time when all human genes (say 100,000 in total) have been found and that simple, diallelic polymorphisms in these genes have been identified. Assume that five such diallelic polymorphisms have been identified within each gene, so that a total of $10 \times 10^5 = 10^6$ alleles need to be tested. The statistical problem is that the large number of tests that need to be made leads to an inflation of the type 1 error probability. For a linkage test with pairs of affected siblings, we use a lod score (logarithm of the odds ratio for linkage) criterion of 3.0, which asymptotically corresponds to a type 1 error probability α of about 10^{-4} . In a linkage genome screen with 500 markers, this significance level gives a probability greater than 95% of no false positives. The equivalent false positive rate for 1,000,000 independent association tests can be obtained with a significance level $\alpha = 5 \times 10^{-8}$.

We illustrate the power of linkage versus association tests at different significance levels by determining the sample size N (number of families) necessary to obtain 80% power (the probability of rejecting the null hypothesis when it is false) (6) (see table). With a linkage approach and a disease gene with a GRR of 4 or greater, the number of affected sibling pairs necessary to detect linkage is realistic (185 or 297), provided the allele frequency p is between 5 and 75%. For a gene with a GRR of 2 or less, however, the sample sizes are generally beyond reach (well

over 2000), precluding their identification by this approach. In contrast, the required sample size for the association test, even allowing for the smaller significance level, is vastly less than for linkage, especially for affected sibling pair families when the value of p is small. Even for a GRR of 1.5, the sample sizes are generally less than 1000, well within reason.

Thus, the primary limitation of genome-wide association tests is not a statistical one but a technological one. A large number of genes (up to 100,000) and polymorphisms (preferentially ones that create alterations in preferred proteins or their expression) must first be identified, and an extremely large number of such polymorphisms will need to be tested. Although testing such a large number of polymorphisms on several hundred, or even a thousand families, might currently seem implausible in scope, more efficient methods of screening a large number of polymorphisms (for example, sample pooling) may be possible. Furthermore, the number of tests we have used as the basis for our calculations (1,000,000) is likely to be far larger than necessary if one allows for linkage disequilibrium, which could substantially reduce the required number of markers and families needed for initial screening.

Some of the important loci for complex diseases will undoubtedly be found by linkage analysis. However, the limitations to detecting many of the remaining genes by linkage studies can be overcome; numerous genetic effects too weak to identify by linkage can be detected by genomic association studies. Fortunately, the samples currently collected for linkage studies (for example, affected pairs of siblings and their parents) can also be used for such association studies. Thus, investigators should preserve their samples for future large-scale testing.

The human genome project can have more than one reward. In addition to sequencing the entire human genome, it can lead to identification of polymorphisms for all the genes in the human genome and the diseases to which they contribute. It is a charge to the molecular technologists to develop the tools to meet this challenge and provide the information necessary to identify the genetic basis of complex human diseases.

References and Notes

1. N. Risch, *Am. J. Hum. Genet.* **40**, 1 (1987); *ibid.* **46**, 229 (1990).
2. From the formulas in (1), we have $\lambda_0 = 1 + 0.5V_A/K^2$ and $\lambda_S = 1 + (0.5V_A + 0.25V_D)/K^2$, where $K = p^2\gamma^2 + 2pq\gamma + q^2 = (p\gamma + q)^2$, $V_A = 2pq(\gamma - 1)^2$ ($p\gamma + q$), and $V_D = p^2q^2(\gamma - 1)^2$. Hence, $\lambda_0 = 1 + w$ and $\lambda_S = (1 + 0.5w)^2$, where $w = pq(\gamma - 1)^2$. The proportion of alleles shared is given by $Y = 1 - 0.5z_1 - z_0$, where z_1 and z_0 are the probabilities of the sib pair sharing 1 and 0 disease alleles ibd, respectively. From (1), $z_0 = 0.25/\lambda_S$ and $z_1 = 0.5\lambda_C/\lambda_S$. Thus, after some algebra, $Y = 1 - 0.25(\lambda_0 + 1)/$

$$\lambda_S = (1 + w)/(2 + w).$$

3. R. Spielman, R. E. McGinnis, W. J. Ewens, *Am. J. Hum. Genet.* **52**, 506 (1993).
4. By Bayes theorem, the probability of a parent of an affected child being heterozygous is given by $P(\text{Het}|\text{Aff child}) = P(\text{Het})P(\text{Aff Child}|\text{Het})/P(\text{Aff Child}) = 2pq(0.5p(\gamma^2 + \gamma) + 0.5q(\gamma + 1))/(p\gamma + q)^2 = pq(\gamma + 1)/(p\gamma + q)$.
5. E. S. Lander and N. J. Schurk, *Science* **265**, 2037 (1994).
6. Consider a set of M independent, identically distributed random variables B_i of discrete value. Under the null hypothesis H_0 , assume $E(B_i) = 0$ and $\text{Var}(B_i) = 1$. Under the alternative hypothesis H_1 , let $E(B_i) = \mu$ and $\text{Var}(B_i) = \sigma^2$. For a sample of size M , let $T = \sum B_i/\sqrt{M}$. Then under H_0 , T also has mean 0 and variance 1, while under H_1 , it has mean $\sqrt{M}\mu$ and variance σ^2 . We assume that T is approximately normally distributed both under H_0 and H_1 . Then the sample size M required to obtain a power of $1 - \beta$ for a significance level α is given by

$$M = (Z_\alpha - \sigma Z_{1-\beta})^2 / \mu^2 \quad (1)$$

For each affected sib pair, we score the number of alleles shared ibd from each of $2N$ parents. Define $B_i = 1$ if an allele is shared from the i th parent and $B_i = -1$ if unshared. Under the null hypothesis of no linkage, $P(B_i = 1) = P(B_i = -1) = 0.5$, so $E(B_i) = 0$ and $\text{Var}(B_i) = 1$. For the genetic model described above with genotypic relative risks of γ^2 , γ , and 1, allele sharing by affected sibs is independent for the two parents; thus, we can consider sharing of alleles one parent at a time. Thus, for affected sib pairs assuming $\theta = 0$ and no linkage disequilibrium, the formula is

$$N = \frac{(Z_\alpha - \sigma Z_{1-\beta})^2}{2\mu^2}$$

where

$$\begin{aligned} \mu &= 2Y - 1 \\ \sigma^2 &= 4Y(1 - Y) \\ Y &= \frac{1 + w}{2 + w} \\ w &= \frac{pq(\gamma - 1)^2}{(p\gamma + q)^2} \end{aligned}$$

$Z_\alpha = 3.72$ (corresponding to $\alpha = 10^{-4}$), and $Z_{1-\beta} = -0.84$ (corresponding to $1 - \beta = 0.80$). For an association test using the transmission/disequilibrium test, with the disease locus or a nearby locus in complete disequilibrium, the number (N) of families with affected singletons required for 80% power is also calculated from formula 1. For this case, we score the number of transmissions of allele A from heterozygous parents. Let h be the probability a parent is heterozygous under the alternative hypothesis, namely, $h = pq(\gamma + 1)/(p\gamma + q)$. Then define $B_i = h^{-0.5}$ if the parent is heterozygous and allele A is transmitted; $B_i = 0$ if the parent is homozygous; and $B_i = -h^{-0.5}$ if the parent is heterozygous and transmits allele a. Under the null hypothesis, $E(B_i) = 0$ and $\text{Var}(B_i) = 1$. Under the alternative hypothesis, $\mu = E(B_i) = \sqrt{h}(\gamma - 1)/(\gamma + 1)$ and $\sigma^2 = \text{Var}(B_i) = 1 - h(\gamma - 1)^2/(\gamma + 1)^2$. In this case, there are two parents per family and they act independently, so the required number (N) of families is given by half of formula 1 where μ and σ^2 are given above. Here, $Z_\alpha = 5.33$ (corresponding to $\alpha = 5 \times 10^{-8}$). For the same test but with affected sib pairs instead of singletons, the number of families required is given by half of formula 1 (transmissions from two parents to two children) with the same formulas for μ and σ^2 as for singleton families but now using the heterozygote frequency for parents of affected sib pairs. Using the above formulas, we can calculate sample sizes for the three study designs.

27 October 1995; accepted 6 June 1996.

The use of a genetic map of biallelic markers in linkage studies

Leonid Kruglyak

BEST AVAILABLE COPY

Improvements in genetic mapping techniques have driven recent progress in human genetics. The use of single nucleotide polymorphisms (SNPs) as biallelic genetic markers offers the promise of rapid, highly automated genotyping. As maps of SNPs and the techniques for genotyping them are being developed, it is important to consider what properties such maps must have in order for them to be useful for linkage studies. I examine how polymorphic and densely spaced biallelic markers need to be for extraction of most of the inheritance information from human pedigrees, and compare maps of biallelics with today's genome-scanning sets of microsatellite markers. I conclude that a map of 700–900 moderately polymorphic biallelic markers is equivalent—and a map of 1,500–3,000 superior—to the current 300–400 microsatellite marker sets.

The revolution in human genetics that has unfolded over the past decade and a half has been driven largely by the development of genetic maps. The original concept was proposed by Botstein *et al.*, with restriction fragment length polymorphisms (RFLPs) as markers¹. The first human RFLP was quickly identified², and Huntington's disease soon became the first autosomal disorder linked to an anonymous DNA marker³. The first RFLP map of the human genome followed shortly⁴. RFLPs were based on a variety of polymorphisms at the sequence level (single nucleotide changes, insertions and deletions, repeat length polymorphisms) and were assayed by Southern hybridization. Although a great advance, RFLPs were often not very polymorphic, and they were costly and time-consuming to develop and assay in large numbers. Nevertheless, these markers made human molecular genetics a reality and led to the mapping of a number of important mendelian diseases.

The next major advance came with the discovery and development of microsatellites (STRs or SSLPs) as markers⁵. These loci are abundant, have fairly high polymorphism rates and can be assayed by PCR, leading to lower cost and a greater degree of automation. Dense maps of microsatellites are now available^{6,7}, allowing simple mendelian diseases to be mapped with relative ease and enabling first searches for genetic causes of complex diseases by genome scan. However, the requirements to assay the loci on gels and to distinguish several length-based alleles make it hard to fully automate the genotyping process, and typing large numbers of individuals for markers covering the genome remains beyond the resources of all but a few labs. There is thus a need to move beyond this current technology.

Recent attention has focused on the use of single nucleotide polymorphisms (SNPs) as genetic markers. At first glance, this may appear to represent a step back to the days of low polymorphism rates characteristic of RFLPs. However, modern technology should allow efficient assays of SNPs in numbers sufficiently large to offset their lower polymorphism rates, as discussed below. SNPs offer a number of important advantages over microsatellites. They are highly abundant, with classic estimates of more than 1 per 1,000 base pairs, or more than 3 million in the genome^{8,9}. To date, more than 1,000 PCR-amplifiable SNP markers have been discovered and mapped (D. Wang, pers. comm.). Because SNPs have only two (common) alleles (hence the term 'biallelics'), genotyping them requires only a plus/minus assay rather than a length measurement, permitting easier automation. Several non-gel-based assays have been proposed^{10–14}, with high-

density oligonucleotide arrays currently showing great promise for typing large numbers of biallelic markers in parallel^{15,16}.

Here I consider the feasibility of carrying out linkage studies with a genetic map based on biallelic markers. The key questions are: What level of polymorphism is required? and How many markers adequately cover the genome? These questions are addressed below.

Assumptions

The effects of marker density and polymorphism were examined by simulating pedigree genotype data and measuring the information content^{17,18} for a broad range of map densities and polymorphism levels (see Methods for simulation details). Information content measures the fraction of inheritance information extracted by the map relative to that which

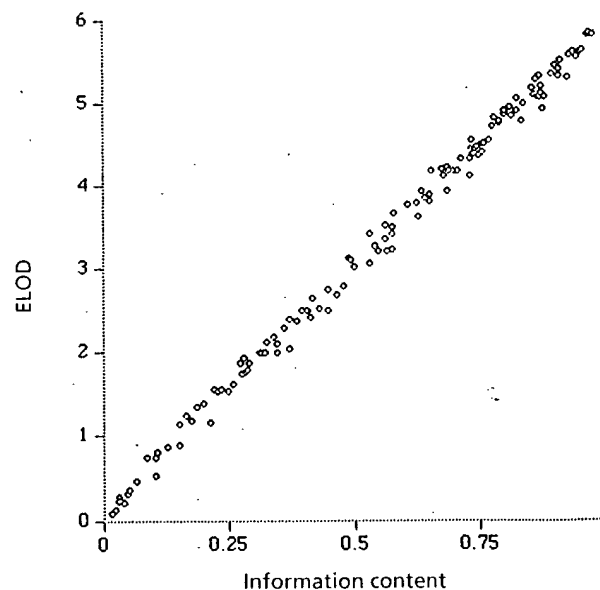


Fig. 1 Expected lod score (ELOD) for a dominant locus is plotted against information content. Each circle represents the results of a simulation for one of 130 maps, as described in Methods. The solid line shows the expected linear correlation if information content of 0 corresponds to an ELOD of 0 and information content of 1 corresponds to the maximum achievable ELOD of 6.02 in these pedigrees.

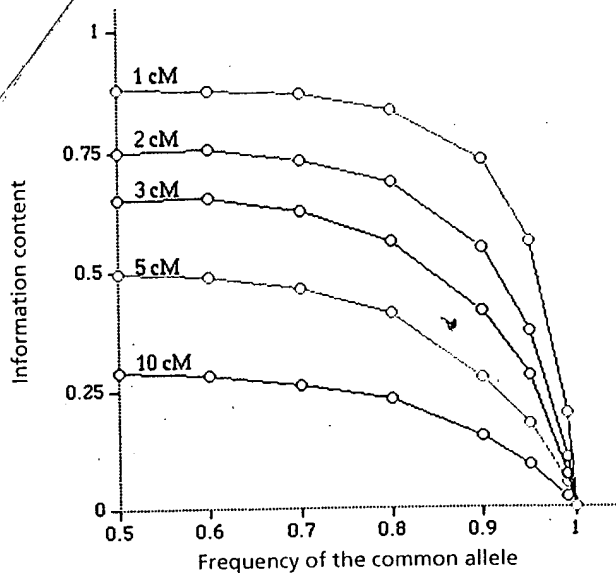


Fig. 2 Information content for five map densities is plotted against the frequency of the more common of the two alleles of a biallelic marker. The circles show actual simulation data points.

would be extracted by an infinitely dense polymorphic map. Thus, an information content of 1 reflects complete information; whereas an information content of 0 reflects no information. Information content incorporates both marker density and polymorphism in a single general measure of map quality that is independent of assumptions about a particular disease locus. It also closely predicts the power of a map to detect linkage—for example, as measured by the expected lod score (ELOD; Fig. 1).

The markers were assumed to be evenly spaced, and information content was measured at a location halfway between two markers, where it is expected to be lowest. For clarity, a single pedigree structure is used throughout: first-cousin pairs with parents but not grandparents available for genotyping. Extensive simulations show that although the absolute numbers differ somewhat for other pedigree structures, all the main conclusions about the relative importance of marker polymorphism and density continue to hold:

How polymorphic do biallelic markers need to be?

Biallelic markers vary in their rates of polymorphism: the more common allele can range in frequency from 50% to nearly 100%. In considering a map of biallelic markers, it is important to ask whether only near-perfect (50-50) biallelics are useful or whether less polymorphic markers can provide comparable amounts of information. To answer this question, I measured information con-

Table 1 • Information content for biallelics

spacing (cM)	allele distribution				
	50-50	60-40	70-30	80-20	90-10
1	0.88	0.88	0.87	0.84	0.73
2	0.75	0.76	0.73	0.69	0.55
3	0.65	0.65	0.63	0.56	0.42
4	0.58	0.56	0.53	0.48	0.34
5	0.50	0.49	0.46	0.41	0.27
6	0.45	0.43	0.41	0.36	0.24
7	0.39	0.39	0.37	0.32	0.22
8	0.35	0.35	0.33	0.28	0.19
9	0.32	0.31	0.29	0.25	0.17
10	0.29	0.28	0.26	0.23	0.15

tent in simulations of maps of biallelic markers with varying degrees of polymorphism.

The results (Fig. 2, Table 1) clearly indicate that at higher map densities, allele frequency has only a small effect on information content in the range of frequency distributions from 50-50 to 80-20. Specifically, a 1-cM map of 60-40 biallelics provides an information content of 0.88, essentially the same as perfect 50-50 biallelics at this density, while 70-30 biallelics provide an information content of 0.87, and 80-20 biallelics provide an information content of 0.84. The information content drops to 0.73 for 90-10 biallelics. Thus, the use of biallelic markers with frequency distribution as skewed as 80-20 leads to little reduction in the information content of a dense map. For sparser maps of 5-10 cM, a similar conclusion holds for marker allele frequency distributions as skewed as 70-30.

How dense does a map of biallelic markers need to be?

Although there is a limit on how polymorphic a biallelic marker can be (a 50-50 distribution of the two alleles), there is essentially no theoretical limit on map density (or marker number), as reasonably polymorphic SNPs can be found roughly every 1 kb, or about 3 million times in the human genome (see above). Thus, one answer to how many markers are needed is that more is always better¹. For common linkage study designs, however, the addition of markers provides diminishing returns once most of the inheritance information has been extracted. As shown above, a 1-cM

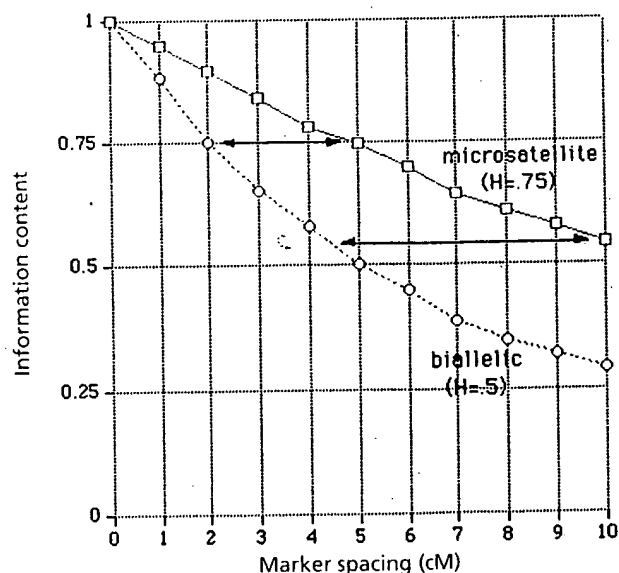


Fig. 3 Information content is plotted against marker spacing for selected microsatellite (heterozygosity $H=0.75$) and biallelic (heterozygosity $H=0.5$) markers. Arrows connect the points on the two curves where information content reaches 0.75 (top) and 0.54 (bottom), the values for 5-cM and 10-cM microsatellite maps, respectively.

Table 2 • Information content for microsatellites

spacing (cM)	allele number			
	3	4	5	10
1	0.93	0.95	0.96	0.97
2	0.87	0.90	0.91	0.94
3	0.80	0.84	0.87	0.90
4	0.74	0.78	0.81	0.87
5	0.68	0.75	0.78	0.82
6	0.64	0.70	0.73	0.80
7	0.58	0.64	0.69	0.76
8	0.53	0.61	0.66	0.72
9	0.49	0.58	0.62	0.69
10	0.45	0.54	0.58	0.68

map of 50–50 biallelic markers extracts 88% of the available information, and it is unlikely that higher information content is needed in an initial screen for linkage. What is the informational cost of decreasing the density of the map? Simulation results (Fig. 2) show that map density plays a more critical role than marker polymorphism. A 2-cM map provides information content of 0.75, a 3-cM map 0.65 and a 5-cM map 0.50. Together with the results of the previous section, these numbers lead to the conclusion that for initial linkage studies it is desirable to screen a dense (1–2-cM) map of moderately polymorphic (50–50 to 80–20) biallelic markers. Interesting regions can then be followed up with all available (biallelic and microsatellite) markers.

It is worth noting that there are two separate issues regarding map density: how many markers exist and how many markers can be genotyped rapidly and cost-effectively. Although current microsatellite maps cover the genome at an average spacing of less than 1 cM (with more than 5,000 markers in the final Génethon map alone⁷), genotyping more than a few hundred markers in a large collection of families remains beyond the power of today's technology and research budgets. Thus, the practical limit on the number of biallelic markers will depend on the techniques for marker development and genotyping. Nonetheless, it is interesting to compare such maps with current maps of microsatellite markers. Such a comparison is carried out in the next section.

Comparison of maps based on biallelics and microsatellites

Current genome scans typically employ a 10-cM map of microsatellite markers for the initial screen^{19,20}, followed by denser coverage of regions that yield interesting results. (Although one could employ a 'staged search' strategy of starting with a sparser 20–40-cM map and then increasing the density in all moderately positive regions^{21,22}, economies of scale in large genotyping labs usually argue for a one-stage initial scan: using a single optimized set of markers for all projects is more efficient than 'filling in' different regions for each.) Microsatellite markers typically vary between 0.65 and 0.8 in heterozygosity (for instance, an average of 0.7 in the final Génethon map⁷), and for simplicity I will use microsatellites with four equally frequent alleles (heterozygosity of 0.75) as representative in the following comparisons with biallelics with two equally frequent alleles (heterozygosity of 0.5); results for other values are given in Tables 1 and 2. Intuitively, one would expect two closely linked biallelics to provide the same information as one microsatellite, and simulations largely confirm this intuition. A 10-cM map of microsatellites achieves information content of 0.54 (Fig. 3). The same information content is provided by a 4.5-cM map of biallelic markers. A denser 5-cM microsatellite map achieves an information content of 0.75, as does a 2-cM map of biallelics. In general, maps of biallelic markers at about 2.25–2.5 times the density of microsatellites provide a comparable information content. A 10-cM map of 300 microsatellite markers can therefore be replaced by a 4-cM map of 750 biallelic markers. These conclusions are in rough agreement with the results of an earlier study of the trade-off between marker spacing and polymorphism²³.

As technology improves, it is likely that screening a much denser map of biallelic markers will be cheaper and easier than carrying out today's genome scans employing microsatellites^{15,16}. There are reasons to employ such denser maps. As shown above, current scan densities lead to considerable loss of information. This problem is more serious for data-sets consisting of more distantly related affecteds or of progeny of consanguineous marriages used in homozygosity mapping²⁴. It is therefore worth noting that a 1-cM map of biallelics (about 3,000 markers) yields much higher information content than a 10-cM map of microsatellites (0.88 vs. 0.54), and is superior to a 5-cM microsatellite map (0.88 vs. 0.75).

Practical linkage analysis using biallelic markers

Because of the lower polymorphism rates of biallelic markers, it is critical to consider many linked markers simultaneously; indeed, all the above results assume complete multipoint analysis of all markers on a chromosome. Such multipoint analysis is even more important for biallelics than for microsatellites. Fortunately, recently developed algorithms and software allow multipoint analysis with an essentially unlimited number of linked markers to be carried out for sib pairs¹⁷ as well as for general pedigrees of moderate size¹⁸. These methods can also be used for automatic haplotype reconstruction, avoiding the tedious prospect of haplotyping many biallelics by hand. The one remaining challenge is extending multipoint analysis with many markers to large multi-generational families, although even here the situation is improving²⁵.

Discussion

The results presented here clearly demonstrate that the use of a genetic map of biallelic markers for linkage studies is feasible on theoretical grounds. It is not necessary to find only 'perfect' 50–50 biallelics: markers with allele frequency distributions as skewed as 70–30 or even 80–20 are almost as useful in a dense map. This result should allay the concern that markers discovered in one population may not be sufficiently informative in other populations with different allele frequencies. A 1–2-cM map of moderately polymorphic biallelic markers is superior to today's microsatellite screening sets for extracting inheritance information and should provide a more efficient tool for initial genome scans.

Even denser maps should enable novel study designs for dissecting genetically complex phenotypes. In particular, genome scans for linkage disequilibrium (LD) and association may become practical^{26–28}. Because LD mapping relies on detecting recombinationally conserved regions around an ancestral mutation, the required map density will vary with the age and history of the study population, with very dense maps (spacing of 10 kb or less) likely to be needed for LD scans in a mixed general population. A more promising approach may be to screen in parallel functional (coding) biallelic polymorphisms in many genes for direct association (rather than LD) with disease^{26–28}.

Maps of biallelic markers and the technology to genotype them should be forthcoming^{15,16}, and the resulting progress in human genetics will be exciting to watch.

Methods

Simulations. Segregation of chromosomes of 100-cM length with evenly spaced markers was simulated. For biallelics, the frequencies of the common allele were 0.5, 0.6, 0.7, 0.8, 0.9, 0.95 and 0.99. For microsatellites, equally frequent alleles were assumed, with allele numbers of 3, 4, 5, 10, 20 and 100. Marker spacings of 1, 2, ..., 10 cM were examined. Each simulation consisted of 100 replicates of 10 cousin pairs each. Information content was computed with GENEHUNTER¹⁸. Information content was measured halfway between the two markers closest to the middle of the chromosome. For ELOD computation, a dominant disease locus with full penetrance, no phenocopies and allele frequency of 0.001 was assumed to lie halfway between two markers, and chromosomes were simulated assuming that both cousins were affected. GENEHUNTER was used to compute multipoint lod scores. The relationship between information content and ELOD is preserved for other assumptions about the disease locus (data not shown). Simulation software used to generate the data is available from the author and can be used to explore additional map properties and pedigree structures.

Acknowledgements

I thank M. Daly, E. Lander and D. Wang for helpful discussions and comments on the manuscript. This work was supported in part by a Special Emphasis Research Career Award from NHGRI (HG00017).

1. Botstein, D., White, D.L., Skolnick, M. & Davis, R.W. Construction of a genetic linkage map in man using restriction fragment length polymorphisms. *Am. J. Hum. Genet.* **32**, 314-331 (1980).
2. Wyman, A.R. & White, R.W. A highly polymorphic locus in human DNA. *Proc. Natl. Acad. Sci. USA* **77**, 6754-6758 (1980).
3. Gusella, J.F., et al. A polymorphic DNA marker genetically linked to Huntington's disease. *Nature* **306**, 234-238 (1983).
4. Donis-Keller, H. et al. A genetic linkage map of the human genome. *Cell* **51**, 319-337 (1987).
5. Weber, J.L. & May, P.E. Abundant class of human DNA polymorphisms which can be typed using the polymerase chain reaction. *Am. J. Hum. Genet.* **44**, 388-396 (1989).
6. Cooperative Human Linkage Center. A comprehensive human linkage map with centimorgan density. *Science* **265**, 2049-2054 (1994).
7. Dib, C. et al. A comprehensive genetic map of the human genome based on 5264 microsatellites. *Nature* **380**, 152-154 (1996).
8. Hofker, M.H. et al. The X chromosome shows less genetic variation at restriction sites than the autosomes. *Am. J. Hum. Genet.* **39**, 438-451 (1986).
9. Cooper, D.N., Smith, B.A., Cooke, H.J., Niemann, J. & Schmidtke, J. An estimate of unique DNA sequence heterozygosity in the human genome. *Hum. Genet.* **69**, 201-205 (1985).
10. Nickerson, D.A. et al. Automated DNA diagnostics using an ELISA-based oligonucleotide ligation assay. *Proc. Natl. Acad. Sci. USA* **87**, 8923-8927 (1990).
11. Livak, K.J., Marmaro, J. & Todd, J.A. Towards fully automated genome-wide polymorphism screening. *Nature Genet.* **9**, 341-342 (1995).
12. Saiki, R.K., Walsh, P.S., Levenson, C.H. & Erlich, H.A. Genetic analysis of amplified DNA with immobilized sequence-specific oligonucleotide probes. *Proc. Natl. Acad. Sci. USA* **86**, 6230-6234 (1989).
13. Syvanen, A.-C., Aalto-Setälä, K., Harju, L., Kontula, K. & Soderlund, H. A primer-guided nucleotide incorporation assay in the genotyping of apolipoprotein E. *Genomics* **8**, 684-692 (1990).
14. Wu, D.Y., Ugozzoli, L., Pal, B.K. & Wallace, R.B. Allele-specific enzymatic amplification of β -globin genomic DNA for diagnosis of sickle cell anemia. *Proc. Natl. Acad. Sci. USA* **86**, 2757-2760 (1989).
15. Wang, D. et al. Toward a third generation genetic map of the human genome based on biallelic polymorphisms. *Am. J. Hum. Genet.* **59**, A3 (1996).
16. Chee, M. et al. Accessing genetic information with high-density DNA arrays. *Science* **274**, 610-614 (1996).
17. Kruglyak, L. & Lander, E.S. Complete multipoint sib-pair analysis of qualitative and quantitative traits. *Am. J. Hum. Genet.* **57**, 439-454 (1995).
18. Kruglyak, L., Daly, M.J., Reeve-Daly, M.P. & Lander, E.S. Parametric and nonparametric linkage analysis: a unified multipoint approach. *Am. J. Hum. Genet.* **58**, 1347-1363 (1996).
19. Reed, P.W. et al. Chromosome-specific microsatellite sets for fluorescence-based, semi-automated genome mapping. *Nature Genet.* **7**, 390-395 (1994).
20. Dubovsky, J., Sheffield, V.C., Duyk, G.M. & Weber, J.L. Sets of short tandem repeat polymorphisms for efficient linkage screening of the human genome. *Hum. Mol. Genet.* **4**, 449-452 (1995).
21. Elston, R.C. Designs for the global search of the human genome by linkage analysis. in *Proceedings of the 16th International Biometrics Conference* 39-51 (Hamilton, New Zealand, 1992).
22. Brown, D.L., Gorin, M.B. & Weeks, D.E. Efficient strategies for genomic searching using the affected-pedigree-member method of linkage analysis. *Am. J. Hum. Genet.* **54**, 544-552 (1994).
23. Terwilliger, J.D., Ding, Y. & Ott, J. On the relative importance of marker heterozygosity and intermarker distance in gene mapping. *Genomics* **13**, 951-956 (1992).
24. Lander, E.S. & Botstein, D. Homozygosity mapping: a way to map human recessive traits with the DNA of inbred children. *Science* **236**, 1567-1570 (1987).
25. O'Connell, J.R. & Weeks, D.E. The VITESSE algorithm for rapid exact multilocus linkage analysis via genotype set-recoding and fuzzy inheritance. *Nature Genet.* **11**, 402-408 (1995).
26. Risch, N. & Merikangas, K. The future of genetic studies of complex human diseases. *Science* **273**, 1516-1517 (1996).
27. Lander, E.S. The new genomics: global views of biology. *Science* **274**, 536-539 (1996).
28. Collins, F.S. Positional cloning moves from perditional to traditional. *Nature Genet.* **9**, 347-350 (1995).

BEST AVAILABLE COPY

D-0.1 M KCl, Tat-SF/pp140 was eluted with increasing salt concentrations and was detected mostly in 0.2 to 0.4 M KCl fractions. These fractions were pooled, dialyzed against buffer D-0.1 M KCl, and loaded onto a glutathione Sepharose (Pharmacia) column containing GST-Tat fusion proteins. After the column was washed with buffer D-0.4 M KCl, Tat-SF/pp140 was eluted from the column with buffer D containing 1.4 M KCl. The estimated overall purification after these steps was ~3000-fold. In the experiment shown in Fig. 3, the 0.2 to 0.4 M KCl heparin Sepharose fraction containing Tat-SF activity was subjected to fractionation through an Affi-Gel 10 matrix column (Bio-Rad) containing immobilized Tat. Tat-SF activity was eluted from the column with increasing salt concentrations. The 0.6 M KCl fraction was analyzed as described in Fig. 3.

10. T. O'Brien, S. Herdin, A. Greenleaf, J. T. Lis, *Nature* 370, 75 (1994); M. E. Dahmus, *Biochim. Biophys. Acta* 1261, 171 (1995).
11. A. P. Rice and F. Carlotti, *J. Virol.* 64, 1864 (1990).
12. The Tat-SF/pp140 fraction eluted from the GST-Tat column was subjected to SDS-polyacrylamide gel electrophoresis (PAGE), and the pp140 polypeptide was blotted onto a nitrocellulose membrane. Approximately 15 µg of pp140 were recovered from the membrane and subjected to digestion with lys-C. Six major peptides were obtained and microsequenced. One of the peptides (KMNAQETATGMAFEERIDE) was contained in the sequence of EST60354 in the Washington University-Merck EST database. An Xho I-Eco RI fragment corresponding to the C-terminus of the Tat-SF1 gene and its 3' untranslated region was labeled and used as a probe to screen a λZidLox (Gibco BRL) cDNA library prepared from human HL60 cells. Complementary DNAs were recovered from seven independent plaques in the autonomously replicating plasmid pZL1 as instructed by the manufacturer (Gibco BRL). The largest cDNA clone containing the full-length Tat-SF1 gene was named pZL-Tat-SF1-4b and was sequenced by dideoxy-DNA sequencing with T7 DNA polymerase.
13. D. R. Marshak and D. Carroll, *Methods Enzymol.* 200, 134 (1991).
14. D. J. Kenan, C. C. Query, J. D. Keene, *Trends Biochem. Sci.* 16, 214 (1991).
15. O. Delattre et al., *Nature* 359, 162 (1992); P. H. Sorensen et al., *Nature Genet.* 6, 146 (1994).
16. A. Crozat, P. Aman, N. Mandahl, D. Ron, *Nature* 363, 640 (1993); T. H. Rabbitts, A. Forster, R. Larson, P. Nathan, *Nature Genet.* 4, 175 (1993).
17. M. Ladanyi, *Diagn. Mol. Pathol.* 4, 162 (1995); T. H. Rabbitts, *Nature* 372, 143 (1994).
18. S. E. Harper, Y. Qiu, P. A. Sharp, *Proc. Natl. Acad. Sci. U.S.A.* 93, 6536 (1996).
19. J. W. Little and M. R. Green, *Nature* 338, 39 (1989).
20. H. Kato et al., *Genes Dev.* 6, 655 (1992); R. A. Marciniak and P. A. Sharp, *EMBO J.* 10, 4189 (1991).
21. M. G. Izban and D. S. Luse, *Genes Dev.* 6, 1342 (1992); D. Wang and D. K. Hawley, *Proc. Natl. Acad. Sci. U.S.A.* 90, 843 (1993).
22. E. Bengal, O. Flores, A. Krauskopf, D. Reinberg, Y. Aloni, *Mol. Cell. Biol.* 11, 1195 (1991); J. Greenblatt, J. R. Nodwell, S. W. Mason, *Nature* 364, 401 (1993).
23. C. H. Herrmann and A. P. Rice, *J. Virol.* 69, 1612 (1995).
24. N. A. McMillan et al., *Virology* 213, 413 (1995).
25. W. A. May et al., *Mol. Cell. Biol.* 13, 7393 (1993); H. Zinsner, R. Albalat, D. Ron, *Genes Dev.* 8, 2513 (1994); D. D. Prasad, M. Ouchida, L. Lee, V. N. Rao, E. S. Reddy, *Oncogene* 9, 3717 (1994).
26. P. J. Mitchell and R. Tjian, *Science* 245, 371 (1989).
27. S. F. Altshul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, *J. Mol. Biol.* 215, 403 (1990).
28. M. A. Truett et al., *RNA* 4, 333 (1985).
29. H. E. Gendelman et al., *Proc. Natl. Acad. Sci. U.S.A.* 83, 9759 (1986).
30. L. S. Tilley, P. H. Brown, B. R. Cutler, *Virology* 178, 560 (1990).
31. J. R. Neumann, C. A. Morency, K. O. Russian, *Bio-Techniques* 5, 444 (1987).
32. We are grateful to B. Pepinsky and Biogen for providing pure HIV Tat protein and Tat mutant TatΔC; to J. Borrow (Massachusetts Institute of Technology (MIT) Center for Cancer Research) for human cDNA libraries; and to R. Cook (MIT Biopolymers Laboratory) for peptide

sequencing. We thank K. Luo, J. Borrow, and H. Kawasaki for valuable advice and discussions; and B. Blencowe, K. Ceppek, G. Jones, K. Luo, and C. Query for helpful comments on the manuscript. We also thank M. Sialaca for secretarial support. Supported by grants from the National Institutes of Health (GM34277 and

AI32486) to P.A.S., and partially supported by a National Cancer Institute Center core grant (CA14051). O.Z. was supported by a postdoctoral fellowship of The Jane Coffin Childs Memorial Fund for Medical Research.

19 June 1996; accepted 23 August 1996

Accessing Genetic Information with High-Density DNA Arrays

Mark Chee, Robert Yang, Earl Hubbell, Anthony Berno, Xiaohua C. Huang, David Stern, Jim Winkler, David J. Lockhart, Macdonald S. Morris, Stephen P. A. Fodor

Rapid access to genetic information is central to the revolution taking place in molecular genetics. The simultaneous analysis of the entire human mitochondrial genome is described here. DNA arrays containing up to 135,000 probes complementary to the 16.6-kilobase human mitochondrial genome were generated by light-directed chemical synthesis. A two-color labeling scheme was developed that allows simultaneous comparison of a polymorphic target to a reference DNA or RNA. Complete hybridization patterns were revealed in a matter of minutes. Sequence polymorphisms were detected with single-base resolution and unprecedented efficiency. The methods described are generic and can be used to address a variety of questions in molecular genetics including gene expression, genetic linkage, and genetic variability.

A central theme in modern genetics is the relation between genetic variability and phenotype. To understand genetic variation and its consequences on biological function, an enormous effort in comparative sequence analysis will need to be carried out. Conventional nucleic acid sequencing technologies make use of analytical separation techniques to resolve sequence at the single nucleotide level (1, 2). However, the effort required increases linearly with the amount of sequence. In contrast, biological systems read, store, and modify genetic information by molecular recognition (3). Because each DNA strand carries with it the capacity to recognize a uniquely complementary sequence through base pairing, the process of recognition, or hybridization, is highly parallel, as every nucleotide in a large sequence can in principle be queried at the same time. Thus, hybridization can be used to efficiently analyze large amounts of nucleotide sequence. In one proposal, sequences are analyzed by hybridization to a set of oligonucleotides representing all possible subsequences (4). A second approach, used here, is hybridization to an array of oligonucleotide probes designed to match specific sequences. In this way the most informative subset of probes is used. Implementation of these concepts relies on recently developed combinatorial technologies to generate any ordered array of a large number of oligonucleotide probes (5).

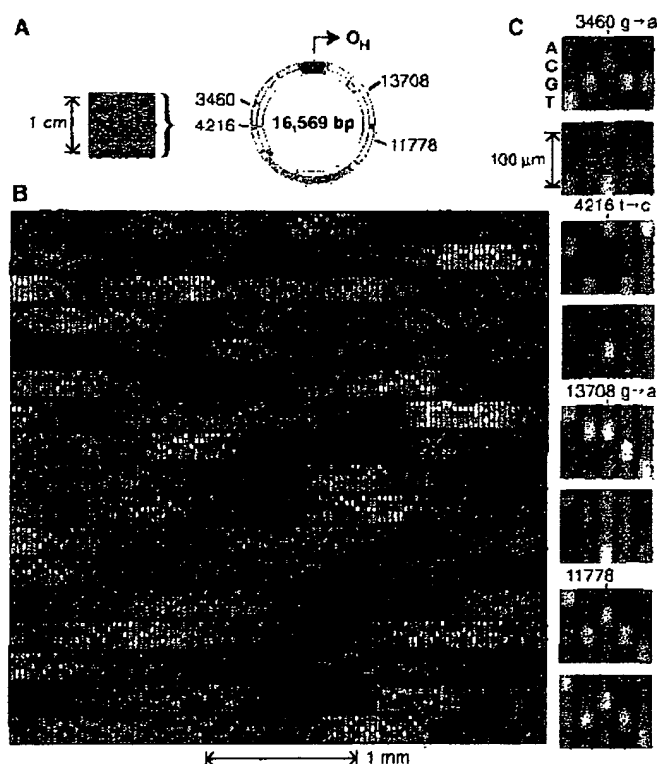
The fundamentals of light-directed oligonucleotide array synthesis have been described (5, 6). Any probe can be synthesized at any discrete, specified location in the array, and any set of probes composed of the four nucleotides can be synthesized in a maximum of $4N$ cycles, where N is the length of the longest probe in the array. For example, the entire set of $\sim 10^{12}$ 20-nucleotide oligomer probes, or any desired subset, can be synthesized in only 80 coupling cycles. The number of different probes that can be synthesized is limited only by the physical size of the array and the achievable lithographic resolution (7).

An array consisting of oligonucleotides complementary to subsequences of a target sequence can be used to determine the identity of a target sequence, measure its amount, and detect differences between the target and a reference sequence. Many different arrays can be designed for these purposes. One such design, termed a 4L tiled array, is depicted in Fig. 1A. In each set of four probes, the perfect complement will hybridize more strongly than mismatched probes. By this approach, a nucleic acid target of length L can be scanned for mutations with a tiled array containing $4L$ probes. For example, to query the 16,569 base pairs (bp) of human mitochondrial DNA (mtDNA), only 66,276 probes of the possible $\sim 10^9$ 15-nucleotide oligomers need to be used.

The use of a tiled array of probes to read a target sequence is illustrated in Fig. 1C. A tiled array of 15-nucleotide oligomers varied

Affymetrix, 3380 Central Expressway, Santa Clara, CA 95051, USA.

Fig. 3. Human mitochondrial genome on a chip. (A) An image of the array hybridized to 16.6 kb of mitochondrial target RNA (L strand). The 16,569-bp map of the genome is shown, and the H strand origin of replication (O_H), located in the control region, is indicated. (B) A portion of the hybridization pattern magnified. In each column there are five probes: A, C, G, T, and Δ , from top to bottom. The Δ probe has a single-base deletion instead of a substitution and hence is 24 instead of 25 bases in length. The scale is indicated by the bar beneath the image. Although there is considerable sequence-dependent intensity variation, most of the array can be read directly. The image was collected at a resolution of ~ 100 pixels per probe cell. (C) The ability of the array to detect and read



single-base differences in a 16.6-kb sample is illustrated. Two different target sequences were hybridized in parallel to different chips. The hybridization patterns are compared for four different positions in the sequence. Only the P^{25,13} probes are shown. The top panel of each pair shows the hybridization of the mt3 target, which matches the chip P⁰ sequence at these positions. The lower panel shows the pattern generated by a sample from a patient with Leber's hereditary optic neuropathy (LHON). Three known pathogenic mutations, LHON3460, LHON4216, and LHON13708, are clearly detected. For comparison, the fourth panel in the set shows a region around position 11,778 that is identical in both samples.

provide the foundation for a powerful genetic analysis technology. The method can be used to characterize the spectrum of sequence variation in a population and can be applied to the analysis of many genes in parallel. In the case of human mtDNA, we simultaneously analyzed the control region, 13 protein coding genes, 22 tRNA genes, and 2 ribosomal RNA genes. The methods described here can be applied to other research areas in molecular genetics; for example, the ability to identify and sequence polymorphisms provides a basis for genetic mapping. The specificity of oligonucleotide hybridization and the scalability of the method suggests the possibility of a dedicated array that could be used to generate a high-resolution genetic map of an entire genome in a single experiment. Likewise, the concepts and techniques described here have been used to develop approaches for mRNA identification and the large-scale, real-time measurement of expression levels (24). Thus, the sequence of a gene, its spectrum of change in the population, its chromosomal location, and its dynam-

ics of expression (all essential to a full understanding of function) can be determined with high-density probe arrays. The challenge now is to synthesize and read probe arrays at even higher density. For example, a 2 cm by 2 cm array, synthesized with probes occupying 1- μ m synthesis sites in a 4L tiling, could query the entire coding content of the human genome, estimated at 100,000 genes.

REFERENCES AND NOTES

1. F. Sanger, S. Nicklen, A. R. Coulson, *Proc. Natl. Acad. Sci. U.S.A.* 74, 5463 (1977).
2. A. M. Maxam and W. Gilbert, *ibid.*, p. 560.
3. J. D. Watson and F. H. C. Crick, *Nature* 171, 737 (1953).
4. W. Bains and G. C. Smith, *J. Theor. Biol.* 135, 303 (1988); Y. P. Lysov et al., *Dokl. Akad. Nauk. SSSR* 303, 1508 (1988); R. Drmanac, I. Labat, I. Brunker, R. Crkvenjakov, *Genomics* 4, 114 (1989); E. Southern, U. Maskos, R. Elder, *ibid.* 13, 1008 (1992); see also R. B. Wallace et al., *Nucleic Acids Res.* 6, 3543 (1979).
5. S. P. A. Fodor et al., *Science* 251, 767 (1991).
6. A. C. Pease et al., *Proc. Natl. Acad. Sci. U.S.A.* 91, 5022 (1994).
7. In the present format, we can routinely achieve a density of 409,600 synthesis sites in a 1.28 cm by 1.28 cm array. Each 20 μ m by 20 μ m site contains

- $\sim 4 \times 10^6$ functional copies of a specific probe, which corresponds to a mean distance of about 100 Å between probes (M. O. Trulsson, D. Stern, R. P. Rava, unpublished results).
8. S. Anderson et al., *Nature* 280, 457 (1981).
9. The control region of mtDNA is characterized by high amounts of sequence polymorphism concentrated in two hypervariable regions [B. D. Greenberg, J. E. Newbold, A. Sugino, *Gene* 21, 33 (1983); C. F. Aquardo and B. D. Greenberg, *Genetics* 103, 287 (1983)].
10. R. L. Cann, W. M. Brown, A. C. Wilson, *Genetics* 106, 479 (1984).
11. The mt1 and mt2 sequences were cloned from amplified genomic DNA extracted from hair roots [P. Gill, A. J. Jeffreys, D. J. Werrett, *Nature* 318, 577 (1985); R. K. Saiki et al., *Science* 239, 487 (1988)]. The clones were sequenced conventionally (7). Cloning was performed only to provide a set of pure reference samples of known sequence. For templates for fluorescent labeling, DNA was reamplified from the clones with primers bearing bacteriophage T3 and T7 RNA polymerase promoter sequences (bold; mtDNA sequences uppercase): L15935-T3, 5'-ctcgaattacccctcactaaaggAAACCTTTTCC-AAGGA and H667-T7, 5'-taatacgcactataggga-gAGGCTAGGACCAACCTATT.
12. Labeled RNAs from the two complementary mtDNA strands [designated L and H (8)] were transcribed in separate reactions from a promoter-tagged polymerase chain reaction (PCR) product. Each 10- μ l reaction contained 1.5 mM each of the triphosphate nucleotides ATP, CTP, GTP, and UTP; 0.24 mM fluorescein-12-CTP (Du Pont); 0.24 mM fluorescein-12-UTP (Boehringer Mannheim); ~ 1 to 5 nM (1.5 μ l) crude unpurified 1.3-kb PCR product; and T3 or T7 RNA polymerase (1 U/ μ l) (Promega) in a reaction buffer supplied with the enzyme. The reaction was carried out at 37°C for 1 to 2 hours. RNA was fragmented to an average size of <100 nucleotides by adjusting the solution to 30 mM MgCl₂, by the addition of 1 M MgCl₂, and heating at 94°C for 40 min. Fragmentation improved the uniformity and specificity of hybridization (M. Chee et al., data not shown). The extent of fragmentation is dependent on the magnesium ion concentration [J. W. Huff, K. S. Sastri, M. P. Gordon, W. E. C. Wacker, *Biochemistry* 3, 501 (1964); J. J. Butzow and G. L. Eichorn, *Biopolymers* 3, 95 (1965)]. Good hybridization results have been obtained with both DNA and RNA targets prepared with a variety of labeling schemes, including incorporation of fluorescent and biotinylated deoxynucleoside triphosphates by DNA polymerases, incorporation of dye-labeled primers during PCR, ligation of labeled oligonucleotides to fragmented RNA, and direct labeling by photo-cross-linking a psoralen derivative of biotin directly to fragmented nucleic acids [L. Wodicka, personal communication].
13. For two-color detection experiments, the reference and unknown samples were labeled with biotin and fluorescein, respectively, in separate transcription reactions. Reactions were carried out as described (12) except that each contained 1.25 mM of ATP, CTP, GTP, and UTP and 0.5 mM fluorescein-12-UTP or 0.25 mM biotin-16-UTP (Boehringer Mannheim). The two reactions were mixed in the ratio 1:5 (v/v) biotin:fluorescein and fragmented (12). Targets were diluted to a final concentration of ~ 100 to 1000 pM in 3M TMACl [W. B. Melchior Jr. and P. H. von Hippel, *Proc. Natl. Acad. Sci. U.S.A.* 70, 298 (1973)], 10 mM Tris-HCl, pH 8.0, 1 mM EDTA, 0.005% Triton X-100, and 0.2 nM control oligonucleotide labeled at the 5' and with fluorescein (5'-CTGAACGGTAG-CATCTTGAC). Samples were denatured at 95°C for 5 min, chilled on ice for 5 min, and equilibrated to 37°C. A volume of 180 μ l of hybridization solution was then added to the flow cell [R. Lipshutz et al., *Biotechniques* 19, 442 (1995)] and the chip incubated at 37°C for 3 hours with rotation at 60 rpm. The chip was washed six times at room temperature with 6 \times SSPE (0.9 M NaCl, 60 mM NaH₂PO₄, 6 mM EDTA, pH 7.4), 0.005% Triton X-100. Phycoerythrin-conjugated streptavidin (2 μ g/ml in 6 \times SSPE, 0.005% Triton X-100) was added and incubation continued at room temperature for 5 min. The chip was washed again



Mailing Certificate for Amendment/Response May 30, 2006 for application
40/037,718

Mailing Certificate

The following documents were sent today, May 30, 2006 by me, Robert O. McGinnis to the USPTO. These documents were placed in an envelope addressed to Mail Stop Amendment, Honorable Commissioner for Patents P.O. Box 1450 Alexandria, VA 22313-1450 with sufficient first class postage for delivery.

The documents are as follows:

- 1) Amendment/Response a total of 16 pages (signed).
- 2) Enclosures total of 17 sheets (2 pages, signed)
- 3) PTO Credit Card Form PTO-2038 with payment of \$350 signed.
- 4) This Mailing Certificate
- 5) Two Return Receipt Postcards

Robert O. McGinnis
Agent of Record
Reg. No. 44, 232
Ph. 406-522-9355